SHARING AND CONNECTING KNOWLEDGE, RATHER THAN (only) DATA

3 October, 2013 TBC Conference

Anthony Brookes: University of Leicester, UK

Today's Healthcare

Medical literature Primary research Clinical experience Pharmacology

Diagnostics



Inconsistent & sub-optimal healthcare

Tomorrow's Healthcare

Medical literature – Primary research – Clinical experience – Pharmacology –

Diagnostics





New and improved capabilities

WHY share data?

- Enables checking & validation of claims
- Enables meta-analysis
- Enables new uses
- Enables data integration
-
- Panton principles: "For science to effectively function, and for society to reap the full benefits from scientific endeavours, it is crucial that science data be made open"

GEN2PHEN Partners (www.gen2phen.org)









GEN2PHEN - Unification and utility of 'Genotype-to-Phenotype' data

BioShaRE

- Science and IT for Biobanking
- Using health data in research

eTOX

EMIF

- **COPD-MAP**
- REQUITE

- Toxicology data and modelling
- Systems medicine KM (COPD)
- Systems medicine KM (cancer)

Research data sharing

Diagnostic lab data sharing

EHR data sharing

Pharma company data sharing

Public sourced data sharing

DATA SHARING IS IMPORTANT BUT DIFFICULT!

Nuffield Council on Bioethics workshop (12 Sept 2013) "Fact finding meeting on privacy and data control issues"

• General:

- IT can provide solutions, but political will to implement complicated by rapid changes in IT
- nothing can offer one hundred per cent protection
- attitudes to data protection vary considerably across different cultures

Data Warehouses:

- centralisation of data has made it easier to illicitly access private information
- frequent data security breaches creates rational public distrust in many data-heavy projects
- there is a significant black market for personal data

• Patients:

- do not necessarily value privacy to the exclusion of all other interests
- it is not clear how well they understand the issues
- are more likely to offer data if given something in return
- Researchers (carrots, sticks, ambition, ability):
 _????????
- Organisations (business model and legal considerations):
 ????????



ENTITY IDENTIFIERS

Data IDs

- The 'Data Object Identifier' (DOI) system, managed by DataCite. Covers a very broad concept of a 'data object' (much more than just traditional publications). Essential for creating the 'web of data'.

Database IDs

- The BioDBCore project by which database IDs can be assigned. Essential if webservices are to start connecting resources effectively.

Human IDs

- The 'Open Researcher Contributor Identifier' (ORCID) system. Launched late 2013, has already issued many tens of thousands of ORCIDs. Soon to be a required author detail when submitting manuscripts. Removes ambiguity over all 'contributors', thereby enabling incentive/reward systems for data sharing, improved knowledge discovery options, and automation of data access control.

Patient IDs

- Various approaches being tested, but none yet close to an ideal, global solution. Needed to resolve identical patient records that relate to different individuals from those that are merely duplicated (recycled) information about a single individual.

BioResource IDs

- Pilot system emerging from GEN2PHEN & BioShaRE, operated by P3G, as a basis for developing BioResource Impact Factor (BRIF) metrics.

Issues that restrict sharing data



- CANNOT: Data owners may not have time nor funding to manually submit data, and/or submission process and requirements too complicated
- WILL NOT: Data owners receive little or no recognition or reward for releasing data, hence little incentive to try
- MUST NOT: Data owners may have good reasons for not sharing data (ethical, legal, competitive edge)

'GWAS Central'

Data Submit Download Help



GWAS CENTRAL 'Publicity Project'

March 2011

570 e-mails sent to authors of 700 studies [30 undeliverable] 19 (3.5%) positive responses (new or corrected data) 9 (1.7%) neutral or negative responses 512 (94.8%) did not reply

January 2012



Improving on the current state-of-the-art



Single Research Team

Data



Consortium Data



5. BIG DATA...



6 @2013 IBM Corporation





Single Research Team

Data

Can/Will all data be digitised, organised, structured... to make them 'sharable' ?





Research Team

Data



Centralised Sharing of Data



Federated Sharing of Data



Realistic Expectation







Openly share the 'existence' rather than the 'substance' of the datathereafter variably manage data access

Connecting Diagnostic Networks

- Need to enable diagnostic laboratories to check whether sequence variants have previously been seen by other labs [with patients with related phenotypes]
- Currently not possible due to difficulties of data sharing between labs, or with others
- Café Variome can solve this...

The Café Variome Solution

- Allows 'open discovery' of the existence (rather than actual substance) of relevant mutation data
- Thereby, enables networks of labs to easily query for the existence of variants, without necessarily revealing the underlying data, thus overcoming issues of patient confidentiality, etc.
- Currently being extended to handle more sophisticated omics/NGS data handling and deep phenotype data

Networks of labs exchanging data



Allows users to check whether the same mutation(s) have previously been seen by other laboratories/patients (with related phenotypes)

Café Variome Central (http://www.cafevariome.org)



Source	Variant Count
Diagnostic laboratories	193
HGMD	57927
UniProt Variants	23558
1000 Genomes	86305
dbSNP (coding)	100103
PhenCode	3694
DMuDB (pilot data)	21
LSDB	57868
FORGE Canada	73
FINDbase	6850
TOTAL:	336592

Café Variome Central is now populated with all major pubic human genome sequence variants, including all public LSDBs

Record Discovery "Menu"

Hor CON Cafe	ne Share Discove	er About≁	Admin Profile Logout
, and the	Var	iant Discovery	
		Select Source:	
	9 All		
BRCA1		Search term:	
Enter a HGNC (e.g. BRC	gene symbol, RefSeq A1 NM_000014.4:c.	Accession/HGVS nomenclature or chrom 4349A>G NM_000014.4 chr12:132344	osomal region 1534534)
		Q Discover Variants	

Data Sharing Models (facilitated, controlled access)



Record Discovery "Menu"

Home Share Disc Cafe Varione	over Abou	ut *			Admin	Profile Logout		
Ve	ariant [Discove	ery					
	Select	Source:	•					
BRCA1	Search	h term:)			
(e.g. BRCA1 NM_000014.4	Enter a HGNC gene symbol, RefSeq Accession/HGVS nomenclature or chromosomal region (e.g. BRCA1 NM_000014.4:c.4349A>G NM_000014.4 chr12:132344534534)							
			J Han Roma Ba					
Source	Open /	Access	Restric	ted Access	Linked	Access		
1000 Genomes Project	0	0	0 😢		0	0		
dbSNP	1401		0	0	0	0		
Diagnostic Variants	0	0	11		0	0		

Administrators Interface



Data Sharing Granularity

ss V	Cafe Varior	ne	Hom	e Share Discover A	About -			Admir	n Profile Logou	rt						
			Shar	ing Policy			×									
	10 10 10		Choo to:	se the sharing policy level	you would like to set y	your selected va	ariants				Discover About	·		Ad	lmin Profile Log	
	Cafe Variome	G	rest	rictedAccess	\$			haring	Actions		ription	Variant Count	Action	Change S	haring Polic	
<u>_</u>	vx1	G				Confirm	Close	enAccess	C m		Project	86306	+ 🤆 🗎	openAccess	restrictedAccess	
1	vx2	F8		c.6857A>G	NM_000132.3	-		openAccess	СШ		iants	478542	+ © ≘	openAccess	restrictedAccess	
1	vx3	MTN	/1	c.205C>T	NM_000252.2			openAccess	СШ		Database	20	+ 0 🕯	openAccess	restrictedAccess	
5	vx4	BRC	A2	c.10202C>T	NM_000059.3			openAccess	C 🛍		isorders Database	1952	+ ଓ 🗎	openAccess	restrictedAccess	
	vx5	BIN	1	c.141C>T	NM_004305.2			openAccess	C mm		atabase	1294	+ 🤇 🗎	openAccess	restrictedAccess	
											onsortium	74	+ 🤇 🗎	openAccess	restrictedAccess	
							hgn	nd	Human Gene	e Mutatio	on Database	57930	+ 0 🗎	openAccess	restrictedAccess	
							lsd	lb	Locus-specifi		Locus-specific Databases		74571	+ 🤅 📋	openAccess	restrictedAccess
							pheno	phencode Phe		PhenCode		91114	+ 🤆 🛢	openAccess	restrictedAccess	
							unip	rot		Uniprot		67197	+ 🤅 🗎	openAccess	restrictedAccess	

Users can control access to variants from individual fine grained level to whole data sets

Café Variome adoption and use

- Requires only three input fields per recorded mutation
 - HGVS name, gene name & reference URL
 - optional phenotype, patient ID, & other fields
 [metadata includes source DB name, curator name, reference transcript]
- Directly installed from a simply installation file
- Supports networks by being installed
 - on a single server
 - multi-site reference data loaded into this menu
 - multi-site reference data kept locally & accessed remotely by web-services
 - on multiple servers
 - multi-site reference data kept locally & searched from these menus
- Data also loaded directly by upload buttons available on
 - Alamut software (Interactive Biosystems)
 - Gensearch (PhenoSystems)

Cafe Variome Network Collaborators

- DMuDB (with Andrew Devereau, Susan Stenhouse)
- Denmark diagnostic network (with Friedrik Wikman)
- German diagnostic network (with Arne Pfeufer)
- French diagnostic network (with Andre Blavier (Interactive Biosoftware))
- Belgium diagnostic network (with Gert Matthijs)
- Canadian diagnostic network (with Forge Canada / Care4Rare)
- Netherlands diagnostic network (with Morris Swertz & Rolf Sijmons)
- Ehlers Danlos Syndrome network (with Raymond Dalgleish)
- Collagen disease network (with Peter Byers)
- Inherited Colon Cancer network (with Finlay Macrae, InSiGHT)

Data Discovery --> Knowledge Sharing

Record Discovery "Menu"

Home Share Disco	over About •	,			Admin	Profile Logout
Va	ariant Di	scover	у			
	Select So	ource:	1			
	Search t	term:				
BRCA1						
Enter a HGNC gene symbol, RefSe (e.g. BRCA1 NM_000014.4	eq Accession/ł c.4349A>G N	HGVS nome NM_000014	nclature or ch .4 chr12:132	nromosomal regi 344534534)	on	
	Q Discover	Variants				or 11/55 (20%)
					16 DA DA	
Source	Open Ac	cess	Restricte	d Access	Linked	Access
1000 Genomes Project	0	0	0	0	0	0
dbSNP	1401		0	8	0	0
Diagnostic Variants	0	0	11		0	0

DataSHIELD: Pooled data analysis without data sharing

• An Analysis Computer (AC) sends iteratively requests for fitting a given model to the Data Computers (DC) on which data are stored



Knowledge Sharing









Data Sharing





Knowledge Sharing



Data Remote Analysis





Data Discovery

DATA	 important, but problematic
SHARING	- prioritise IDs and risk categorisation

- 90% of the benefits DATA DISCOVERY
 - 10% of the problems
- just a technical extension KNOWLEDGE - 99% of the benefits SHARING
- One does not have to co-locate:
- Data repository & Data interpretation
- One does not have to co-locate:
- Data sharing & Knowledge sharing & Data discovery



Partners in GEN2PHEN, BioShaRE, EMIF, eTox, COPD-MAP, REQUITE

Anthony Brookes' Bioinformatics Group

Tim Beck, Charalambos Chrysostomou, Robert Hastings, Owen Lancaster, Adam Webb Sirisha Gollapudi, Robert Free, Gudmundur Thorisson

DKP Facility members



Data to Knowledge for Practice (DKP) Facility